# Deploying Natural Language Processing for Social Science Analysis

Karin Verspoor, Los Alamos National Laboratory

Antonio Sanfilippo, Pacific Northwest National Laboratory

Mark Elmore, Oak Ridge National Laboratory

Ed Mackerrow, Los Alamos National Laboratory

# Content Analysis for communication research

- Content analysis:
    - methodology for systematic analysis of characteristics of messages
    - used by social scientists to investigate the nature of communications in terms of quantitative variables from which inferences can be drawn about context, meaning, and/or intents
    - Examples (from Neuendorf, 2002):
        - Prevalence of violence in top-grossing films of the 1990s
        - Importance of physical attractiveness in personal ads in newspapers
        - Analysis of the language of schizophrenics

- Harold Lasswell on communication:

*Who* says *what*, to *whom*, *why*, to *what extent*, and with *what effect*?

# Methods for content analysis

- Fundamentally quantitative in nature
  - Neuendorf: "A content analysis has as its goals a numerically based summary of a chosen message set."
  - counts of words or phrases in a given dictionary/category set
  - measurements of the amounts of variables, as determined from *codes* annotated on *messages*
    - This can be meta-data, e.g. date, duration, role or age of a character in a story, socio-economic status of a speaker, religious affiliation
    - Some of these attributes may be directly identifiable in the message

- CAQDAS: Computer-Assisted Qualitative Data Analysis Software
  - Software tools that support dictionary-based and manual coding
  - Frequency output
  - Searching for boolean co-occurrences of codes
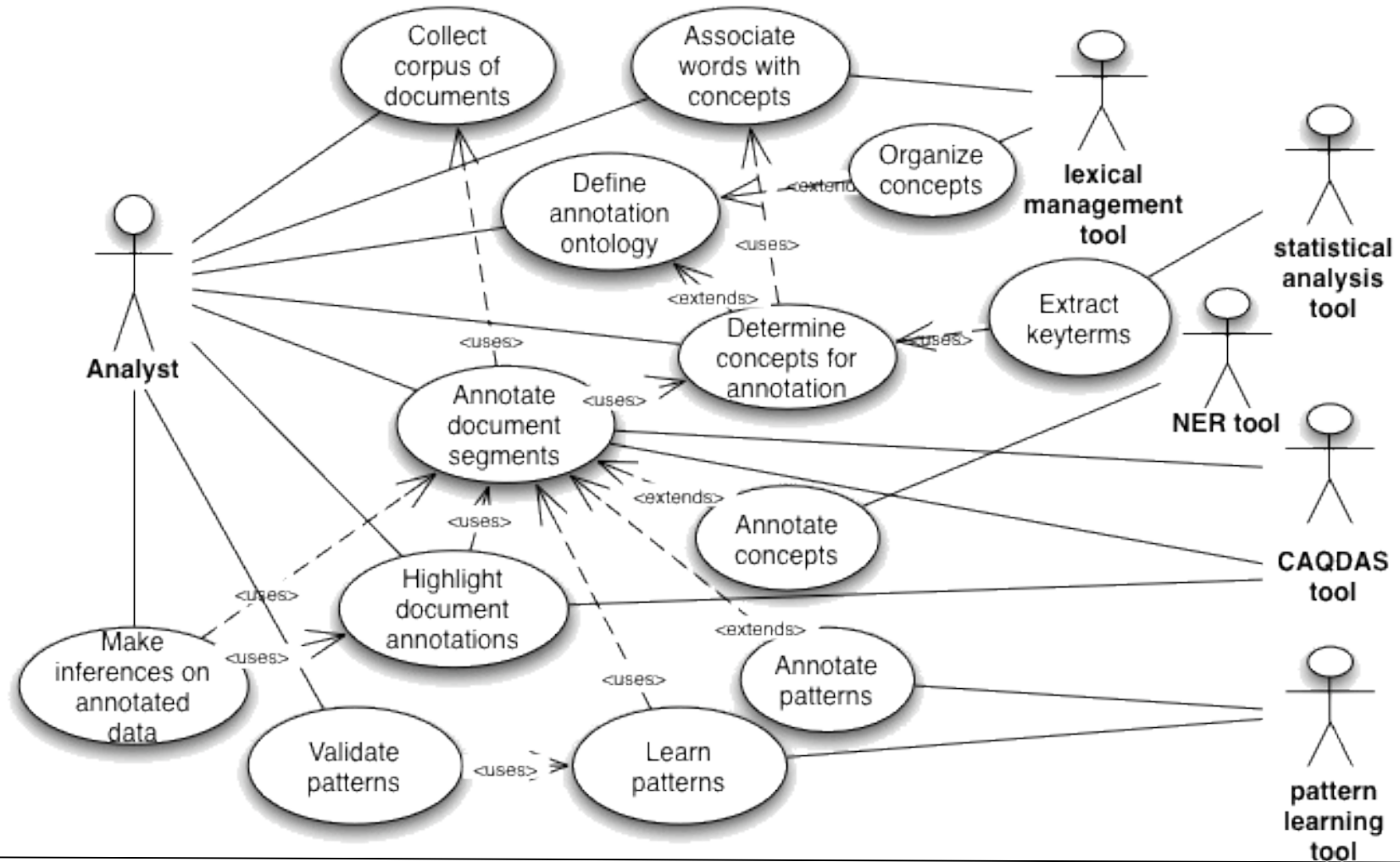  - Much coding, in particular of more complex patterns, must still occur manually

# Natural Language Processing for content analysis

- Natural Language Processing (NLP) aims to exploit linguistic analysis techniques for text to determine
  - *syntactic* (structural) properties
  - *semantic* (meaning) properties
  - *statistical* properties

- Augmentation of content analysis tools with NLP techniques can facilitate higher-throughput content analysis and ultimately more automation of content analysis
  - named entity recognition
  - word-sense disambiguation

- Insights from the NLP research community can further enhance computer-assisted content analysis
  - organizing concepts (lexical management, ontologies)
  - learning information extraction patterns

# Use case for NLP-augmented content analysis

# Augment CAQDAS with NLP

- Integrate NLP tools into the CAQDAS environment
  - Retain manual coding and boolean search capabilities: final decisions rest with the content analyst
  - Bring in more sophisticated text analysis tools to
    - facilitate message vocabulary assessment
    - automatically code many core semantic categories

- Utilize coded messages as a training set for further automation
  - generalize from examples using linguistic properties
  - identify patterns that correspond to more complex (e.g. relational) codes

# Lexical management for content analysis

- Beyond dictionaries: hierarchically organized semantic networks of words
  - Synonyms
  - Hyper/hyponyms

- for concept refinement or generalization

- enabling inference over codes

- manage surface variants of terms "behind the scenes"
  - singular vs. plural
  - verb inflections

- Application-specific annotation ontology
  - Provides a controlled vocabulary
  - Ensures greater consistency in coding

# Statistical analysis for keyterm extraction

- Highly frequent words in a text are not always meaningful (cf. function words, "stop" words)

- Statistical tests can help to hone in on the real content-bearing terms, both words and phrases
  - TF.IDF for words
  - Chi-squared test of independence of words in a bigram

Associate words with concepts

Define annotation ontology

Extract keyterms

statistical analysis tool

<uses>

<extends>

<uses>

Determine concepts for annotation

Analyst

- Statistically derived keyterms can form the foundation of the annotation ontology for the topic of the content analysis

# Named Entity Recognition (NER): *Who* says What to *Whom*?

- NER tools aim to *automatically* identify and label certain basic semantic categories of terms in natural language text
  - People
  - Organizations
  - Places
  - Things (e.g. diseases, weapons, etc.)



- These will likely be useful categories to start coding from

- Allows the computational system to make a first pass at coding; freeing the human coder to identify more complex codes

- Using supervised learning and sufficient training data, named entity recognizers can be trained to recognize other categories of terms relevant to the content analysis domain

# Pattern Learning



- Use annotations in coded data subset as seed patterns

- Use bootstrap semi-supervised learning algorithm to find best matches with seed patterns in large data set
  - Thelen & Riloff, 2002
  - Using linguistic context of codes

# Case Study: Extraction of Social Movement Theory Signatures

- Annotation study undertaken at Pacific Northwest National Laboratory
  - *Antonio Sanfilippo, Andrew Cowell, Annie Boek and Stephen Tratz*

- Identify and code relevant text segments in harvested documents
  - Create ground truth data set semi-manually
  - Utilize NLP techniques to facilitate manual coding
  - Validate manual coding through inter-coder agreement

- Construct queries to find signatures of Social Movement Theory (e.g. Political Constraints and Frames) and quantify their strength

| | |
|---|---|
| • **Political constraint**<br>  – **Source** = Egyptian Regime<br>  – **Policy** = Repression<br>  – **Target** = Sinai Bedouins<br><br>• **Strength** = Supported by 20% of the data consulted<br><br>• **Provenance** = [list document sources] | • **Diagnostic anti-system Frame**<br>  – **Promoter** = Muslim Brotherhood<br>  – **Intent** = Denounce<br>  – **Issue** = Repression<br>  – **Target** = Egyptian Regime<br>  – **Audience** = Egyptian Citizens<br><br>• **Strength =** Supported by 75% of the data consulted<br><br>• **Provenance** = [list document sources] |

# Case Study: Methodology

- Selected 1000 documents about recent terrorist attacks in Egypt

- Developed scheme of 53 codes with several hundred code keywords using insights from Social Movement Theory

- Assessed coding scheme using the Cohen kappa test to measure inter-coder agreement

- Hand-annotated 36 documents

- Formulated searches to identify and quantify frame signatures

- Utilizing the CAQDAS tool Qualrus

- Testing inter-annotator agreement
  - Three people coded the same text (931 words)
  - Used Cohen Kappa test to measure agreement in the assignment of codes to the same segments (Target Kappa Coefficient ≥ 0.7)

| | Cohen Kappa Test Results |
|---|---|
| Coders A & B | 0.78 |
| Coders A & C | 0.69 |
| Coders B & C | 0.73 |

*PNNL researchers: Antonio Sanfilippo, Andrew Cowell, Annie Boek and Stephen Tratz*

# Case Study: Harvest, classify and select documents

# Case Study: Create coding guidelines

- Each code is associated with keywords (concepts in WordNet lexical semantic graph)

- Weights specify the similarity of the keyword to the code

# Case Study: Annotate Documents

- Named Entity Recognition helps identify automatically segments to be coded

- Code keywords and weights derived from lexical management aid manual annotation

# Case Study: Construct queries

What evidence have we that the Egyptian regime is persecuting the Bedouins?

# Case Study: Frame Annotation

| | |
|---|---|
| *The Parliamentary Bloc of the Muslim Brotherhood (MB)* denounces *the insistence of the* security apparatus *on terrorizing innocent people and on using the emergency law against* honest Egyptian citizens, *through its campaign of raids and detentions against Muslim Brothers in the governorates of Cairo, Alexandria, Daqahliya and lastly Minya.* | **PROMOTER** |
| | **INTENTION** **OBJECTIVE** |
| | **TARGET** (the person or organization to blame for grievances) |
| | **ISSUE** |
| | **AUDIENCE** (the group at which the message is aimed) |

- **Diagnostic anti-system Frame**
  - **Promoter** = Muslim Brotherhood
  - **Intent** = Denounce
  - **Issue** = Repression
  - **Target** = Egyptian Regime
  - **Audience** = Egyptian Citizens
- **Strength =** Supported by 75% of the data consulted
- **Provenance** = [list document sources]

- Goal: annotate larger constructs of Social Movement Theory, specifically *Frames*

- Frame constituents
  - **promoter:** *who is behind the message?*
  - **target:** *who is to blame for grievances?*
  - **intention:** *what is the objective of the message?*
  - **issue:** *what are the grievances?*
  - **audience**: *aggrieved group*

- Reason/generalize over language to establish constituent codes (e.g. *issue: repression*)

# Case Study: Assess Frame Evidence

# Conclusions

- NLP techniques dovetail with existing methods in content analysis

- By augmenting CAQDAS tools with NLP, coding can proceed more rapidly
  - automation of basic coding
  - sophisticated tools for handling term identification and variations
  - yet, the human is still in the loop for verification

- Construction of a repository of coded texts for a specific domain can be drawn on to enable learning for further coding automation